

EMPIRICAL ANALYSIS OF ROUGH SET CATEGORICAL CLUSTERING  
TECHNIQUES BASED ON ROUGH PURITY AND VALUE SET

JAMAL UDDIN

A thesis submitted in  
fulfillment of the requirement for the award of the  
Doctor of Philosophy in Information Technology

Faculty of Computer Science and Information Technology  
Universiti Tun Hussein Onn Malaysia

AUGUST 2017

I would like to dedicate my this PhD thesis to my beloved parents whose sincere prayers make it possible for me to be a successful computer and mathematics researcher. May Allah always bless them.



## ACKNOWLEDGEMENT

Surely All praise is for Allah Almighty Who is the creator of this universe and Darood and Salam upon the Holy prophet Hazrat Muhammad (PBUH). Thanks to Allah Almighty Who enabled me to complete this research thesis with the continuous guidance and sincere cooperation of my kind supervisor Prof. Madya Dr Rozaida Binti Ghazali and my co-supervisor Prof. Dr. Mustafa Bin Mat Deris. In fact I learnt a lot from my worthy supervisors whose valuable suggestions, constructive comments, thought provoking ideas gave me this chance to make my research work a success with in a stipulated period. They are real role model and mentor for me.

I am also thankful to Allah Almighty Who bestowed upon me highly talented and sincere teachers, whose selfless guidance and tireless effort gave me opportunities and chances to complete this valuable research work. I extend my heartiest gratitude to my teachers especially the dean of faculty Prof. Madya Dr. Nazri Bin Mohd Nawi and Prof. Madya Dr Rathiah Binti Hashim for their sincere support, valuable comments and encouraging attitude. I am thankful to all my worthy teachers from the core of my heart they all remained cooperative honestly throughout my research work.

I will never forget the educational facilities and research oriented environment provided by the Faculty of Computer Science and Information Technology (FSKTM) and the Universiti Tun Hussien Onn Malaysia (UTHM). The sincere and continuous efforts of UTHM staff and administration to make available all modern and latest facilities to impart quality education in all fields are remarkable. It was their sincere efforts and approach that has made us able to learn information technology (IT) research and complete the research work under the guidance of able IT researchers, who are renowned in Malaysia and outside of the country in their field.

At last, it is a little sad to say that I could not get time for my old parents to serve them, but it is a fact that their prayers always remain with me and are still there with me throughout my education and research work. Their sincere prayers and advice helped me in every step of my life. Their prayers console me very much and I dedicate

this research work to them. I am really thankful to my sisters, brother, wife and kids whose best wishes and prayers always remain with me. I would also like to pay my best regard to my research fellow Rashid Naseem, my colleagues and all friends whose best wishes and prayers always remained with me throughout my research work. I would also like to thank Dr. Muhammad Imran who develop the uthmthesis  $\LaTeX$  project for making the thesis writing process a lot easier for me.



## ABSTRACT

Clustering a set of objects into homogeneous groups is a fundamental operation in data mining. Recently, attention has been put on categorical data clustering, where data objects are made up of non-numerical attributes. The implementation of several existing categorical clustering techniques is challenging as some are unable to handle uncertainty and others have stability issues. In the process of dealing with categorical data and handling uncertainty, the rough set theory has become well-established mechanism in a wide variety of applications including databases. The recent techniques such as Information-Theoretic Dependency Roughness (ITDR), Maximum Dependency Attribute (MDA) and Maximum Significance Attribute (MSA) outperformed their predecessor approaches like Bi-Clustering (BC), Total Roughness (TR), Min-Min Roughness (MMR), and standard-deviation roughness (SDR). This work explores the limitations and issues of ITDR, MDA and MSA techniques on data sets where these techniques fails to select or faces difficulty in selecting their best clustering attribute. Accordingly, two alternative techniques named Rough Purity Approach (RPA) and Maximum Value Attribute (MVA) are proposed. The novelty of both proposed approaches is that, the RPA presents a new uncertainty definition based on purity of rough relational data base whereas, the MVA unlike other rough set theory techniques uses the domain knowledge such as value set combined with number of clusters (NoC). To show the significance, mathematical and theoretical basis for proposed approaches, several propositions are illustrated. Moreover, the recent rough categorical techniques like MDA, MSA, ITDR and classical clustering technique like simple K-mean are used for comparison and the results are presented in tabular and graphical forms. For experiments, data sets from previously utilized research cases, a real supply base management (SBM) data set and UCI repository are utilized. The results reveal significant improvement by proposed techniques for categorical clustering in terms of purity (21%), entropy (9%), accuracy (16%), rough accuracy (11%), iterations (99%) and time (93%).

## ABSTRAK

Pengelompokan satu set objek ke dalam kumpulan homogen adalah operasi asas dalam perlombongan data. Kebelakangan ini, perhatian banyak diberikan kepada pengelompokan data berasaskan kategori, iaitu objek data terdiri daripada atribut bukan berangka. Kebanyakan pelaksanaan teknik pengelompokan berasaskan kategori sedia ada adalah mencabar kerana sebahagiannya tidak dapat mengendalikan isu-isu ketidakpastian dan mempunyai masalah kestabilan. Dalam proses berurusan dengan data berasaskan kategori dan pengendalian ketidakpastian, teori set kasar telah menjadi mekanisme yang mantap dalam pelbagai aplikasi termasuk pangkalan data. Kategori set kasar berdasarkan teknik pengelompokan data seperti Teori-Informatik Bersandarkan Kekasaran (ITDR), Atribut Bersandarkan Maksimum (MDA) dan Atribut Signifikan Maksimum (MSA) telah mengatasi teknik-teknik terdahulu seperti Dwi-Kelompok (BC), Kekasaran Mutlak (TR), Kekasaran Min-Min (MMR), dan Kekasaran Sisihan-Piawai (SDR). Kajian ini membentangkan kekangan dan isu-isu bagi teknik-teknik ITDR, MDA dan MSA ke atas set data tertentu di mana teknik-teknik ini gagal untuk memilih atau menghadapi kesukaran untuk memilih kelompok atribut yang terbaik. Selanjutnya, dua teknik alternatif yang dinamakan Pendekatan Ketulenan Kasar (RPA) dan Atribut Nilai Maksima (MVA) bagi mengelompokkan data berasaskan kategori telah dicadangkan. Pembaharuan bagi kedua-dua teknik yang telah dicadangkan ini adalah berikut; mencadangkan definisi ketidakpastian baharu berdasarkan ketulenan bagi kekasaran pangkalan data hubungan, manakala MVA berbeza dengan teknik teori set kasar lain, yang mana teknik ini menggunakan pengetahuan domain seperti set nilai yang bergabung dengan beberapa kelompok (NoC) dalam memilih kelompok atribut yang terbaik. Bagi menunjukkan signifikasinya, asas matematik dan teori bagi pendekatan yang dicadangkan, beberapa cadangan telah digambarkan. Selain itu, teknik-teknik berasaskan kategori yang terkini seperti MDA, MSA, ITDR dan teknik pengelompokan klasik seperti K-mean asas digunakan sebagai perbandingan dan keputusan perbandingan dibentangkan di dalam

bentuk jadual dan grafik. Bagi kegunaan eksperimen, set data daripada kajian-kajian terdahulu digunakan seperti Supply Base Management (SBM) dan pangkalan data UC Irvine Machine Learning Repository (UCI). Keputusan menunjukkan prestasi bagi teknik yang dicadangkan adalah lebih baik dalam memilih atribut kelompok dan mengelompokkan data berasaskan kategori dari segi ketulenan (21%), entropi (9%), lelaran (99%), masa (93%), ketepatan (16%), dan ketepatan kekasaran (11%).



## CONTENTS

<b>DECLARATION</b>	ii
<b>DEDICATION</b>	iii
<b>ACKNOWLEDGEMENT</b>	iv
<b>ABSTRACT</b>	vi
<b>ABSTRAK</b>	vii
<b>CONTENTS</b>	ix
<b>LIST OF TABLES</b>	xii
<b>LIST OF FIGURES</b>	xiv
<b>LIST OF SYMBOLS AND ABBREVIATIONS</b>	xv
<b>LIST OF PUBLICATIONS</b>	xvi
<b>CHAPTER 1 INTRODUCTION</b>	<b>1</b>
1.1 Research background	1
1.2 Research Motivation	5
1.3 Research Objectives	6
1.4 Research Scope	7
1.5 Research Significance	7
1.6 Thesis Organization	8
<b>CHAPTER 2 LITERATURE REVIEW</b>	<b>9</b>
2.1 Introduction	9
2.2 Cluster analysis	9
2.2.1 Probabilistic and Generative Models	10
2.2.2 Distance-Based Algorithms	11
2.2.3 Density and Grid-Based Methods	13
2.2.4 Software Model Clustering	13
2.2.5 Matrix Factorization and Co-Clustering	14
2.2.6 Related work on cluster analysis	14



2.3	Supplier base management (SBM)	21
2.4	Cluster validation	22
2.4.1	Unsupervised measures	23
2.4.2	Supervised measures	23
2.4.3	Relative measures	23
2.4.4	Related work on cluster validation	24
2.4.5	Accuracy	27
2.4.6	Entropy	28
2.4.7	Purity	28
2.4.8	Rough accuracy	29
2.4.9	Minimum iterations, respond time and Big O notation	29
2.5	Rough set theory	30
2.5.1	Information system	31
2.5.2	Indiscernibility relation	33
2.5.3	Set approximations	34
2.5.4	Related work on rough set theory	37
2.6	Rough categorical data clustering and related work	43
2.7	Comparative analysis of existing rough set categori- cal clustering techniques	47
2.7.1	The ITDR technique	49
2.7.2	Research questions on MDA and MSA techniques	50
2.8	Discussion: Scenario Leading to the Research Framework	63
2.9	Summary	68

## **CHAPTER 3 RESEARCH METHODOLOGY 69**

3.1	Introduction	69
3.2	Proposed research framework	69
3.3	Comparison of rough set based proposed and existing techniques	71
3.4	Information-theoretic purity measure and Rough Purity Approach (RPA)	73
3.5	Maximum value attribute (MVA)	79
3.6	Summary	87

## **CHAPTER 4 EXPERIMENTAL RESULTS AND DISCUSSIONS 88**

4.1	Introduction	88
-----	--------------	----

4.2	Experimental setup	88
4.3	Experiments with rough purity approach (RPA)	89
4.3.1	Computational complexity in terms of time and iterations for RPA technique	90
4.3.2	Other evaluation measures for RPA technique	91
4.4	Experiments with maximum value attribute (MVA)	95
4.4.1	Effect of NoC on the purity and entropy of clustering	95
4.4.2	Small cases for MVA technique	100
4.4.3	Real and UCI data sets for MVA technique	107
4.4.4	Computational complexity in terms of time and iterations for MVA technique	109
4.4.5	Other evaluation measures for MVA technique	114
4.4.6	Comparison with Simple K Mean	115
4.5	Summary	117
<b>CHAPTER 5 CONCLUSION</b>		<b>119</b>
5.1	Accomplished objectives	120
5.1.1	Objective 1	120
5.1.2	Objective 2	121
5.1.3	Objective 3	121
5.2	Contributions of research	122
5.3	Threats to validity	124
5.4	Future Works	124
<b>REFERENCES</b>		<b>126</b>
<b>Vita</b>		<b>140</b>

## LIST OF TABLES

2.1	Summary of related work on cluster analysis	20
2.2	An information system	32
2.3	Information System of Flu Patients	33
2.4	Summary of related work on rough set theory	42
2.5	Summary of existing work on rough categorical data clustering	48
2.6	A Dengue Diagnosis Data Set	54
2.7	Dependency Degree of Attributes for Dengue Diagnosis Information System	55
2.8	Significance Degree of Attributes for Dengue Diagnosis Information System	56
2.9	Dependency Degree of Attributes for Lenses Data Set	57
2.10	Significance Degree of Attributes for Lenses Data Set	57
2.11	Suraj's LEMS Data Set	58
2.12	The attribute dependency degrees from Suraj's LEMS Data Set	59
2.13	The degree of significance of all attributes from Suraj's LEMS Data Set	59
2.14	Grzymala's Information System	60
2.15	The attribute dependency degrees from Grzymala's Information System	60
2.16	Pawlak's Car performance Data Set	60
2.17	The attribute significance degree from Pawlak's Car performance information system	61
2.18	Stores Characterization Data Set	62
2.19	Minimum Iterations and Respond Time for Store Data Set	62
2.20	Minimum Iterations and Time for Train Data Set	63
2.21	Strengths and limitations of rough categorical clustering techniques	65
3.1	Comparison of proposed and existing rough set based techniques	74
3.2	Student's enrollment qualification information system	76
3.3	MMP roughness of Table 3.2	78

3.4	Flu patients data set	85
3.5	Value Set Cardinality of Table 3.4	86
3.6	Comparison of RPA and MVA techniques on Balloons data set	87
4.1	The discretized supplier data set	90
4.2	Time complexity of all techniques	91
4.3	Iterative complexity of all techniques	92
4.4	Suraj's Flu Patients Data (Suraj, 2004)	96
4.5	Stores Characterization Data Set (Pawlak, 1991)	97
4.6	Grzymala Data Set (Grzymala-busse, 2005)	98
4.7	Pawlak's Modified Data Set (Pawlak et al., 1995)	99
4.8	Pawlak's Modified Data Set	100
4.9	Dependency Degree of Attributes from Table 4.8	101
4.10	Significance Degree of Attributes from Table 4.8	101
4.11	Cardinality of Value Sets from Table 4.8	101
4.12	Influenza Data Set	102
4.13	Dependency Degree of Attributes from Table 4.12	103
4.14	Significance Degree of Attributes from Table 4.12	104
4.15	Cardinality of Attributes Value Sets from 4.12	104
4.16	Toys attitude data set	105
4.17	Dependency Degree of Attributes from Table 4.16	106
4.18	Significance Degree of Attributes from Table 4.16	106
4.19	Cardinality of Attributes Value Sets from Table 4.16	107
4.20	Number of Clusters Obtained	109
4.21	Computational complexity comparison of techniques	112
4.22	Iterative complexity of all data sets	113
4.23	Time complexity of all data sets	114
4.24	Comparative performance of techniques for all data sets	116

## LIST OF FIGURES

2.1	A rough set	35
2.2	The ITDR algorithm	50
2.3	The MDA algorithm	51
2.4	The MSA algorithm	52
2.5	Stores characterization data set evaluations	62
2.6	Train Data Set Evaluations	63
2.7	Scenario Leading to the Research Framework	67
3.1	Detail of research process	70
3.2	Proposed research framework	72
3.3	Algorithmic steps comparison	73
3.4	The RPA algorithm	77
3.5	The MVA algorithm	80
4.1	The accuracy of MDA, MSA, ITDR and RPA	93
4.2	The entropy of MDA, MSA, ITDR and RPA	93
4.3	The purity of MDA, MSA, ITDR and RPA	93
4.4	The rough accuracy of MDA, MSA, ITDR and RPA	94
4.5	Evaluation Performance of Surajs Flu Dataset	96
4.6	Evaluation Performance of Stores Characterization Data Set	97
4.7	Evaluation Performance of Grzymala's Data Set	98
4.8	Evaluation Performance of Pawlak's Modified Data Set	99
4.9	Pawlak's Modified Data Set Evaluation Graphs	102
4.10	Pawlak's Influenza Data Set Evaluation Graphs	105
4.11	Infant Toy Attitude Data Set Evaluation Graphs	108
4.12	Clusters visualization of balloons data set	109
4.13	Clusters visualization of soya been data set	110
4.14	Clusters visualization of lenses data set	110
4.15	Clusters visualization of balance scale data set	111
4.16	Comparison of MVA with K-mean technique	117

## LIST OF SYMBOLS AND ABBREVIATIONS

RST	–	Rough Set Theory
MDA	–	Maximum Dependency Attribute
MSA	–	Maximum Significance Attribute
ITDR	–	Information Theory Dependency Roughness
RPA	–	Rough Purity Approach
MVA	–	Maximum Value Attribute
NoC	–	Number of Clusters
BC	–	Bi-Clustering
MMR	–	MinMin-Roughness
SDR	–	Standard Deviation Roughness
U	–	Universe of objects
SBM	–	Supply Base Management



PT TA UTHM  
PERPUSTAKAAN TUNKU TUN AMINAH

## LIST OF PUBLICATIONS

1. Jamal Uddin, Rozaida Ghazali, Mustafa Bin Mat Deris (2016), An Empirical Analysis of Rough Set Categorical Clustering Techniques , *PLoS ONE*, Accepted, DOI: 10.1371/journal.pone.0164803 (ISI Q1, IF=3.54)
2. Jamal Uddin, Rozaida Ghazali, Mustafa Bin Mat Deris, Rashid Naseem, Habib Shah (2016), A Survey on Bug Prioritization, *Artificial Intelligence Review*, Springer, PP 1-36, DOI: 10.1007/s10462-016-9478-6 (ISI Q2, IF=2.1)
3. Jamal Uddin, Rozaida Ghazali, Mustafa Bin Mat Deris, Tutut Herawan (2016), Does Number of Clusters Affects the Purity and Entropy of Clustering?, *International Conference on Soft Computing and Data Mining (SCDM)*, Vol 549, ISBN : 978-3-319-51279-2, Springer Conference.



PTTA UTHM  
PERPUSTAKAAN TUNKU TUN AMINAH

## **CHAPTER 1**

### **INTRODUCTION**

#### **1.1 Research background**

In this present information age, it is believed that information prompts success and strength. The modern technologies like computers and satellites are collecting tremendous amounts of information for us. However, these huge amounts of data in disparate structures overwhelmed in recent years rapidly. Therefore, data base management system (DMBS) and organized data bases are developed (Zaïane, 1999). An efficient DMBS contributes towards effective retrieval of specific information from huge corpus of data. Dealing with huge collections of data, the needs such as automatic summarization of data, discovery of patterns in raw data and extraction of information helps in making better managerial choices. Different kinds of information are collected daily that includes scientific data, software engineering data, games, personal data, satellite sensing, digital media, text reports, business transactions, medical data, world wide web repositories, virtual worlds, surveillance video and pictures.

This enormous amount of data stored in databases, files and other repositories requires a powerful means for interpretation of such data, analysis and for the knowledge extraction that could help in decision-making. Knowledge discovery in databases (KDD) refer to the extraction of previously unknown but potentially useful information which is nontrivial and implicit from the data in databases (Zaïane, 1999). The data mining term being part of the knowledge discovery process is frequently used as synonyms for KDD. The KDD process includes steps like raw data collections leading to formation of new knowledge, data cleaning, data integration, data selection,



data transformation, data mining, pattern evaluation and knowledge representation. The data mining task that is employed determines the kind of information needed to be discovered. In general, there are two types of data mining tasks, that is descriptive and predictive tasks (Zaïane, 1999). Descriptive data mining tasks describe the general properties of the existing data, while the predictive data mining tasks attempts to make predictions based on inference on available data.

Many issues are still pending to be addressed like security, social, interface, mining methodologies, performance and data source before the data mining develops into a conventional and trusted discipline (Zaïane, 1999). The data mining functionalities include prediction, association analysis, classification, clustering, characterization and discrimination etc. The clustering is actually used to analyze accurately the data generated by different modern sources and has appeared as a powerful meta-learning tool. It is considered to be a concise model of the data in the absence of specific labeled information. In particular, the key objective of clustering is to categorize data into clusters so that similar objects are grouped in the same cluster according to specific metrics (Fahad *et al.*, 2014). The internal homogeneity and the external separation is considered by most researchers while describing a cluster (Xu & Wunsch, 2005; Norušis, 2011) i.e., similar objects in the same cluster while different objects in separate clusters.

The different clustering techniques can be broadly classified into partitioning, hierarchical, density, grid and model based approaches (Fahad *et al.*, 2014). Partitioning-based techniques specify the initial groups by reallocating them towards a union and all clusters are determined promptly. In hierarchy based clustering, depending on the medium of proximity the data is organized in a hierarchical manner. Similarly, density-based based approaches separates the data objects based on their regions of density, boundary and connectivity. Grid based technique divides the space of the data objects into grids. Whereas, in model based clustering techniques the fit between the given data and some (predefined) mathematical model is optimized (Fahad *et al.*, 2014). Many domains like academic result analysis of institutions, machine learning, image mining, medical dataset, software engineering, bioinformatics,

information retrieval and pattern recognition uses the core methodology of clustering (Wong *et al.*, 2000; Dharmarajan & Velmurugan, 2013; Naseem *et al.*, 2013; Britto *et al.*, 2014; Aggarwal & Reddy, 2014).

The particular choice of a clustering technique also relies tremendously on specific data type. The different data types are textual, discrete sequences, time series, uncertain data, categorical and multimedia data (Aggarwal & Reddy, 2014). There are several clustering techniques developed to combine objects of same characteristics, however the implementation of them is challenging due to certain issues like categorical data clustering, handling uncertainty, stability and efficiency issues. Different techniques for clustering data having only numerical values were proposed by Haimov *et al.* (1989); Wong *et al.* (2000); Shuanhu *et al.* (2004). Unlike numerical data, the multi-valued attributes known as categorical data have common values or common objects and association between both. To deal with categorical data, a number of clustering techniques have been developed (Huang, 1998; Guha, S.; Rastogi, 1999; Ganti & Ramakrishnan, 1999; Gibson & Kleinberg, 2000). Though, they contributed well to clustering process but they are not able to handle uncertainty (Herawan *et al.*, 2010a). In many cases where there is no sharp boundary between clusters, the uncertainty becomes an important real world issue.

Huang, Gupta and Kang (Huang, 1998; Kim *et al.*, 2004) explored fuzzy sets to handle uncertainty in categorical data clustering. However, to attain the stability and to control the membership fuzziness these techniques require multiple runs (Herawan *et al.*, 2010a). Zdzislaw Pawlak introduced rough set theory (RST) (Pawlak, 1991; Pawlak *et al.*, 1995), a mathematical tool to deal with vagueness and uncertainty. Many researchers and practitioners are attracted towards RST by contributing essentially to the applications and development in the fields of artificial intelligence, decision support systems, machine learning, knowledge acquisition, decision analysis, pattern recognition, expert systems, cognitive sciences, inductive reasoning, and knowledge discovery from data bases (Pawlak & Skowron, 2007). Many interesting applications, the basic ideas of RST and its extensions can be found in several books, issues of the transactions on rough sets, special issues of other journals, international conferences,

proceedings and tutorials (Pawlak & Skowron, 2007).

The RST is a viable system to deal with uncertainty in clustering process of categorical data. RST was originally a symbolic data analysis tool now being developed for cluster analysis (Düntsche & Gediga, 2015). In rough categorical clustering, mainly the data set is expressed as the decision table by introducing a decision attribute. Most of these methods assume one or more given partitions of the data set aiming to find a cluster which best represents the data according to some predefined measure. Set approximation and reduct based methods are the two main ideas of the rough set model which are promising for applications. Tolerance rough set clustering (Ngo & Nguyen, 2004) and rough-K-Means clustering (Peters, 2006) are the examples of set approximation methods. Despite of having satisfactory results, these methods have issues as they depend on several parameters and thresholds (Düntsche & Gediga, 2015). The reduct based methods either work as pre-processing tool or as a tool for cluster generation but the problem of time complexity has not been solved yet (Düntsche & Gediga, 2015).

In RST, a subset of universe can be represented in terms of equivalence classes as clustering of universe. Therefore, RST has been successfully applied for selecting best suitable clustering attribute. The pioneer techniques to select clustering attribute are developed by Mazlack *et al.* (2000) which includes Total Roughness (TR) and Bi-Clustering (BC). These techniques work on the accuracy of roughness (approximation accuracy average) in the RST. Later on, another rough categorical clustering approach named Min-Min Roughness (MMR) was proposed by Parmar *et al.* to improve previous techniques (Parmar *et al.*, 2007). Despite of MMR's better performance, issues like accuracy, computational complexity and purity are yet to be addressed.

In 2010, a technique based on the dependency of attributes was introduced by Herawan *et al.* (2010a) named maximum dependency of attributes (MDA) which uses rough set information system for categorical data clustering. Hassanein and Elmelegy in 2013, proposed maximum significance of attributes (MSA) that utilized the RST concept of significance of attributes for selecting clustering attribute (Hassanein & Elmelegy, 2013). Recently, Park and Choi introduced information-theoretic

dependency roughness (ITDR) technique (Park & Choi, 2015b) which finds the entropy roughness to choose the suitable clustering attribute. It is another rough clustering technique that uses the information-theoretic dependencies of categorical attributes in information systems.

## 1.2 Research Motivation

Today the world is full of data and every day people encounter a large amount of information and they store or represent it as data for further analysis and management. One of the vital means in dealing with these data is to classify or group them into a set of categories or clusters. Rough Set Theory (RST) is a powerful mathematical tool proposed by Pawlak (Pawlak & Skowron, 2007) successfully applied to deal with vagueness and uncertainty in data analysis. The concept of rough set theory in this research work is utilized in terms of data in an information system.

Rough set theory has the ability of decision making in the presence of uncertainty and vagueness. Moreover, it can represent a subset of universe in terms of equivalence classes of partition of the universe. Obviously, every subset of attributes induces unique indiscernibility relation which is an equivalence relation and hence, induces unique clustering. This notion of indiscernibility is very attractive, since each indiscernible relation is also a sort of cluster. In this study, the indiscernibility is used as a measure of similarity without any distance function for clustering the objects.

Recently, the problem of clustering categorical data has received much attention in many fields from statistics to psychology. The categorical data unlike numerical data cannot be naturally ordered. Therefore, those clustering techniques dealing with numerical data cannot be used to cluster categorical data. In addition, very less work has been done for clustering the categorical data. A well-known approach for clustering categorical data is using rough set theory (Park & Choi, 2015a). Originally the motivation and inspiration for this study came from exploring useful limitations and issues of existing rough categorical clustering techniques (Mazlack *et al.*, 2000; Parmar *et al.*, 2007; Herawan *et al.*, 2010a; Hassanein & Elmelegy, 2013; Park & Choi,

2015b). This research is conducted in order to come with more general, efficient and better rough categorical clustering techniques. The MDA, MSA and ITDR techniques outperformed their previous techniques such as BC, TR, MMR etc, however, they have certain issues like accuracy, purity, generalizability and computational complexity. On several data sets, these techniques fail or face difficulties in choosing the suitable clustering attribute. Some of the limitations are outlined:

1. MDA technique cannot perform well on data sets with attributes having zero or equal dependency value.
2. MSA technique also fails to select clustering attribute on data sets having attributes with zero or equal significance value.
3. ITDR techniques face issues like random attribute selection and integrity of classes due to presence of entropy measure.

Accordingly in this work, two rough set based categorical clustering techniques are proposed. The first one, information theoretic Rough Purity Approach (RPA) is introduced by establishing a new rough set metric of uncertainty which is rough purity for categorical data clustering. The proposed RPA technique relates the concept of information theoretic purity to rough sets. Considering the domain knowledge of the data set, the second technique Maximum Value Attribute (MVA) is proposed. Here, the rough value set of an attribute is combined with number of clusters. This technique chooses the suitable clustering attribute on basis of maximum number of clusters by an attribute. Several propositions and experiments on benchmark data sets demonstrate the significance, novelty and contribution of these proposed techniques to practical systems.

### 1.3 Research Objectives

The research objectives are listed as follows:

1. To propose a new rough set based categorical clustering technique Rough Purity Approach that takes into account the purity of attributes.
2. To propose another rough set based categorical clustering technique Maximum

Value Attribute that takes into account the value set of attributes combined with number of clusters.

3. To elaborate the performance of proposed techniques on real and benchmark datasets by comparing them with the recent baseline rough categorical clustering techniques like Maximum Dependency Attribute, Maximum Significant Attribute and Information Theoretic Dependency Roughness and classical K-mean clustering algorithm using accuracy, purity, rough accuracy, time and iterative complexity (Big O notation) and entropy.

#### **1.4 Research Scope**

This research only focuses on proposing two rough set theory based categorical clustering techniques named RPA and MVA. The proposed and existing MDA, MSA, ITDR and classical K-mean techniques are analyzed on several benchmark (UCI and KEEL repositories) and a real Supply Base Management (SBM) data set. The experimental results are evaluated using metrics like accuracy, purity, rough accuracy, number of iterations, respond time and entropy.

#### **1.5 Research Significance**

The system implementation is significant by two ways in this research. Firstly, information-theoretic purity is introduced as a new definition to measure the uncertainty using RST for categorical data clustering. Secondly, a domain knowledge about data like rough value set is utilized to develop another rough categorical clustering technique combined with internal evaluation measure like number of clusters. Both these approaches show significant improvement for clustering categorical data not only in terms of time and iterations but also in terms of accuracy, purity, entropy and rough accuracy.



## 1.6 Thesis Organization

The remaining thesis is arranged as below:

Chapter 2 discusses some fundamental concepts and overview of existing work on clustering the categorical data using RST. It comprises of an information system notion in rough relational database, an indiscernibility relation, set approximations and quality of approximations. This chapter discusses the literature review of existing work for cluster analysis, cluster validation, SBM, RST and rough categorical data clustering. Moreover, it also presents the analysis and limitations of some existing rough categorical data clustering techniques with the help of examples.

Chapter 3 discusses the proposed techniques of clustering the categorical data, named Rough Purity Approach (RPA) and Maximum Value Attribute (MVA). The notion of purity using rough set theory and the value set cardinality are presented. Moreover, the evaluation metrics used in this research are also defined. Several propositions and examples are illustrated to show the significance of proposed techniques.

Chapter 4 illustrates the results of experiments on proposed techniques. An empirical study on ten small, fifteen benchmark data sets and a real SBM data set demonstrates the better performance of proposed techniques. Moreover, they are compared with most recent and leading rough set-based categorical clustering techniques. All the experimental results are discussed and analyzed in detail by presenting them in form of tables and graphs.

Finally, Chapter 5 gives concluding remarks, accomplished objectives, contributions and future work.

## **CHAPTER 2**

### **LITERATURE REVIEW**

#### **2.1 Introduction**

The previous chapter demonstrated that cluster analysis and rough set clustering techniques are widely utilized for numerical and categorical data in various forms. Accordingly, this chapter gives an overview of related work on cluster analysis, validation criteria, rough set theory and rough categorical data clustering.

This chapter comprises of nine sections. An overview of cluster analysis techniques and existing work on them are discussed in Section 2.2. The existing work in the field of supply base management is summarized in Section 2.3. The detail of cluster evaluation measures are described in Section 2.4. Similarly, Section 2.5 explains some preliminaries and related research work on rough set theory. Section 2.6 illustrates the overview of existing research on categorical data clustering. Section 2.7 presents the analysis of best recent rough categorical clustering techniques to explore their limitations. Section 2.8 discusses the scenario that leads to research framework. Section 2.9 summarizes the chapter.

#### **2.2 Cluster analysis**

Clustering is one of the most important unsupervised learning tasks in which the objects are divided into clusters so that similar objects are combined in the same cluster while dissimilar objects in separate clusters. Clustering is widely used in many fields, such as text mining (Naresh Kumar Nagwani, 2012), image analysis (Li



## REFERENCES

- Abawajy, J., Kelarev, A., Chowdhury, M. U., & Herbert, F. J. (2016). Enhancing predictive accuracy of cardiac autonomic neuropathy using blood biochemistry features and iterative multitier ensembles. *IEEE Journal of Biomedical and Health Informatics*, 20(1), 408–415.
- Abawajy, J. H., Kelarev, A., & Chowdhury, M. (2014). Large iterative multitier ensemble classifiers for security of big data. *IEEE Transactions on Emerging Topics in Computing*, 2(3), 352–363.
- Abawajy, J. H., Kelarev, A. V., & Chowdhury, M. (2015). Multistage approach for clustering and classification of ECG data. *Computer Methods and Programs in Biomedicine*, 112(3), 720–730.
- Aggarwal, C., & Reddy, C. (2014). *Data Clustering: Algorithms and Applications*. CRC Press Taylor & Francis Group.
- Aggarwal, C., & Yu, P. (2009). A Survey of Uncertain Data Algorithms and Applications. *IEEE Transactions on Knowledge and Data Engineering*, 21(5), 609–623.
- Agresti, A. (2007). *An introduction to categorical Data Analysis*, vol. 2.
- Ahmad, A., & Dey, L. (2007). A k-mean clustering algorithm for mixed numeric and categorical data. *Data and Knowledge Engineering*, 63(2), 503–527.
- Ahn, Y.-Y., Han, S., Kwak, H., Moon, S., & Jeong, H. (2007). Analysis of Topological Characteristics of Huge Online Social Networking Services. In *Proceedings of the 16th International Conference on World Wide Web*. pp. 835–844.
- Aldana-Bobadilla, E., & Kuri-Morales, A. (2015). A clustering method based on the maximum entropy principle. *Entropy*, 17(1), 151–180.
- Amigó, E., Gonzalo, J., Artiles, J., & Verdejo, F. (2009). A comparison of extrinsic

- clustering evaluation metrics based on formal constraints. *Information Retrieval*, 12(4), 461–486.
- Anaraki, J., & Eftekhari, M. (2013). Rough set based feature selection: A Review. In *Proceedings of 5th Conference on Information and Knowledge Technology (IKT)*. pp. 301–306.
- Anquetil, N., & Lethbridge, T. C. (1999). Experiments with Clustering as a Software Remodularization Method. In *Proc. Sixth Working Conf. Reverse Eng.*, pp. 235–255.
- Astel, A., Tsakovski, S., Barbieri, P., & Simeonov, V. (2007). Comparison of self-organizing maps classification approach with cluster and principal components analysis for large environmental data sets. *Water Research*, 41(19), 4566–4578.
- Babcock, B., Datar, M., Motwani, R., & O’Callaghan, L. (2003). Maintaining variance and k-medians over data stream windows. *Proceedings of the twenty second ACM symposium on Principles of database systems*, pp. 234–243.
- Banitaan, S. (2013). TRAM : An Approach for Assigning Bug Reports using their Metadata. pp. 215–219.
- Bean, C., & Kambhampati, C. (2008). Autonomous clustering using rough set theory. *International Journal of Automation and Computing*, 5(1), 90–102.
- Beaubouef, T., Petry, F. E., & Arora, G. (1998). Information-theoretic measures of uncertainty for rough sets and rough relational databases. *Journal of Information Sciences*, 5.
- Berry, M. W. (2004). *Survey of Text Mining : Clustering, Classification, and Retrieval*.
- Biernacki, C., Celeux, G., & Govaert, G. (2006). Assessing a mixture model for clustering with the integrated completed likelihood. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(7), 719–725.
- Bozcan, O., & Bener, A. B. (2013). Handling missing attributes using matrix factorization. In *Proceedings of 2nd International Workshop on Realizing Artificial Intelligence Synergies in Software Engineering (RAISE)*. IEEE. pp. 49–55.

- Britto, D., Filho, F., Carvalho, E., Alexandre, J., Paranhos, R., Batista, M., Sofia, B., & Duarte, F. (2014). Cluster Analysis for Political Scientists. *Applied Mathematics*, 5(August), 2408–2415.
- Cameron, a. C., Gelbach, J. B., & Miller, D. L. (2006). Bootstrap-Based Improvements for Inference with Clustered Errors. *Review of Economics and Statistics*, 90(3), 414–427.
- Cardot, H., Cénac, P., & Monnez, J.-M. (2012). A fast and recursive algorithm for clustering large datasets with k-medians. *Computational Statistics and Data Analysis*, Volume 56,(Issue 6), 1434–1449.
- Chang, H. T., & Peng, H. W. (2012). Facial Image Prediction Using Exemplar-based Algorithm and Non-negative Matrix Factorization. In *Signal & Information Processing Association Annual Summit and Conference (APSIPA ASC)*. pp. 1–4.
- Chen, G., Jaradat, S., Banerjee, N., Tanaka, T., Ko, M., & Zhang, M. (2002). Evaluation and comparison of clustering algorithms in analyzing es cell gene expression data. *STATISTICA SINICA*, 12(1), 241–262.
- Chen, L.-F., & Tsai, C.-T. (2016). Data mining framework based on rough set theory to improve location selection decisions: A case study of a restaurant chain. *Tourism Management*, 53, 197–206.
- Chowdhury, M., Abawajy, J., Kelarev, A., & Jelinek, H. (2016). A Clustering-Based Multi-Layer Distributed Ensemble for Neurological Diagnostics in Cloud Services. *IEEE Transactions on Cloud Computing*, 4(2), 1–1.
- Christopher D. Manning, P. R., & Schütze, H. (2009). *Introduction to Information Retrieval*.
- Darshit et al., P. (2010). A clustering algorithm for supplier base management. *International Journal of Production Research*, 48(13), 3803–3821.
- Davey, J., & Burd, E. (2000). Evaluating the suitability of data clustering for software remodularisation. In *Proceedings of Seventh Working Conference on Reverse Engineering*. IEEE Comput. Soc. pp. 268–276.
- Deris, M. M., Abdullah, Z., Mamat, R., & Yuan, Y. (2015). A new limited tolerance

- relation for attribute selection in incomplete information systems. In *2015 12th International Conference on Fuzzy Systems and Knowledge Discovery (FSKD)*. IEEE. pp. 964–970.
- Dharmarajan, A., & Velmurugan, T. (2013). Applications of partition based clustering algorithms: A survey. *Proceedings of IEEE International Conference on Computational Intelligence and Computing Research*.
- Düntsche, I., & Gediga, G. (2015). Rough set clustering. Tech. rep., Brock University Department of Computer Science Rough, Ontario, Canada.
- Fahad, A., Alshatari, N., Tari, Z., Alamri, A., Khalil, I., Zomoya, A., Foufou, S., & Bauras, A. (2014). A Survey of Clustering Algorithms for Big Data: Taxonomy & Empirical Analysis. *IEEE Transactions on Emerging Topics in Computing*, 2(3), 1–13.
- Feng, H., Chen, Y., Ni, Q., & Huang, J. (2014). A New Rough Set Based Classification Rule Generation Algorithm (RGI). In *Proceedings of International Conference on Computational Science and Computational Intelligence*. Ieee. pp. 380–385.
- Feng, J., & Seok, H. (2011). Applying agglomerative hierarchical clustering algorithms to component identification for legacy systems. *Information and Software Technology*, 53(6), 601–614.
- Ganti, V., & Ramakrishnan, J. G. R. (1999). CACTUS Clustering Categorical Data Using Summaries. In *Proceedings of the fifth ACM SIGKDD international conference on Knowledge discovery and data mining*. pp. 73–83.
- Gao, Y., Zhang, X., Wu, L., Yin, S., & Lu, J. (2017). Resource basis, ecosystem and growth of grain family farm in China: Based on rough set theory and hierarchical linear model. *Agricultural Systems*, 154(May), 157–167.
- Garcia, H. V., & Shihab, E. (2014). Characterizing and Predicting Blocking Bugs in Open Source Projects Categories and Subject Descriptors. In *Proceedings of the 11th Working Conference on Mining Software Repositories*. pp. 72–81.
- Gibson, D., & Kleinberg, J. (2000). Clustering categorical data : an approach based on dynamical systems. *The VLDB Journal*, 8, 222–236.

- Gong, X., & Zhang, G. (2016). Non-Negative Matrix Co-Factorization for Weakly Supervised Image Parsing. *IEEE Signal Processing Letters*, 9908(c), 1–1.
- Grzymala-busse, J. W. (2005). Rough Set Theory with Applications to Data Mining. *Real World Applications of Computational Intelligence, Volume 179*, pp 221–244.
- Guha, S., Meyerson, A., Mishra, N., Motwani, R., & OCallaghan, L. (2003). Clustering data streams: Theory and practice. *Knowledge and Data Engineering, IEEE Transactions on*, 15(3), 515–528.
- Guha, S., Mishra, N., Motwani, R., & O'Callaghan, L. (2000). Clustering data streams. In *Proceedings 41st Annual Symposium on Foundations of Computer Science*. IEEE Comput. Soc. pp. 359–366.
- Guha, S.; Rastogi, R. K. S. (1999). ROCK : A Robust Clustering Algorithm for Categorical. In *Proceedings., 15th International Conference on Data Engineering.*, pp. 512 – 521.
- Haimov, S., Michalev, M. A., & Savchenko, A. (1989). Classification of radar signatures by autoregressive model fitting and cluster analysis. *IEEE Transactions on Geoscience and Remote Sensing*, 21(5), 606–610.
- Har-Peled, S., & Mazumdar, S. (2004). Coresets for k-Means and k-Median Clustering and their Applications. *Proceedings of the thirty-sixth annual ACM symposium on Theory of computing*, pp. 291—300.
- Hassanein, W., & Elmelegy, A. (2013). An Algorithm for Selecting Clustering Attribute using Significance of Attributes. *International Journal of Database Theory & Application*, 6(5), 53–66.
- Herawan, T., Deris, M. M., & Abawajy, J. H. (2010a). A rough set approach for selecting clustering attribute. *Knowledge-Based Systems*, 23(3), 220–231.
- Herawan, T., Ghazali, R., Tri, I., Yanto, R., & Deris, M. M. (2010b). Rough Set Approach for Categorical Data Clustering 1. *International Journal of database theory and Application*, 3(1), 179–186.
- Herawan, T., Tri, I., Yanto, R., & Deris, M. M. A. T. (2010c). ROSMAN : ROugh Set approach for clustering Supplier base MANagement. *Biomedical Soft*

*Computing and Human Sciences*, 16(2), 105–114.

Hodge, V. J., & Austin, J. I. M. (2004). A Survey of Outlier Detection Methodologies. *Artificial Intelligence Review*, 22, 85–126.

Huang, A. (2008). Similarity measures for text document clustering. *Proceedings of the Sixth New Zealand Computer Science Research Student Conference*, (April), 49–56.

Huang, Z. (1998). Extensions to the k -Means Algorithm for Clustering Large Data Sets with Categorical Values. *Data Mining and Knowledge Discovery*, 2, 283–304.

Hunter, M. G., & Peters, G. (2012). Rough Sets: Selected Methods and Applications in Management and Engineering. *Advanced Information and Knowledge Processing*, pp. 129–138.

Jain, A. K. (2010). Data clustering: 50 years beyond K-means. *Pattern Recognition Letters*, 31(8), 651–666.

Jia, X., Shang, L., Zhou, B., & Yao, Y. (2016). Generalized attribute reduct in rough set theory. *Knowledge-Based Systems*, 91, 204–218.

Jyoti (2013). Clustering categorical data using rough set: A Review. *International Journal of Advanced Research in IT and Engineering*, 2(12), 30–37.

Karaboga, D., & Ozturk, C. (2011). A novel clustering approach: Artificial Bee Colony (ABC) algorithm. *Applied Soft Computing*, 11(1), 652–657.

Kent (2008). Cluster Validation. Tech. rep.

Khatami, A., Mirghasemi, S., Khosravi, A., & Nahavandi, S. (2015). A New Color Space Based on K-Medoids Clustering for Fire Detection. *Proceedings of IEEE International Conference on Systems, Man, and Cybernetics, SMC 2015*, pp. 2755–2760.

Kim, D.-w., Lee, K. H., & Lee, D. (2004). Fuzzy clustering of categorical data using fuzzy centroids. *Pattern Recognition Letters*, 25(11), 1263–1271.

Komorowski, J., Polkowski, L., & Skowron, A. (1999). Rough sets: A tutorial. *Rough fuzzy*, pp. 2–8.



- Kontostathis, A., Galitsky, L. M., Pottenger, W. M., Roy, S., & Phelps, D. J. (2004). *A survey of emerging trend detection in textual data mining*.
- Krause, D. R., Handfield, R. B., & Scannell, T. V. (1998). An empirical investigation of supplier development: reactive and strategic processes. *Journal of Operations Management*, 17(1), 39–58.
- Kumar, P., & Tripathy, B. (2009). MMeR an algorithm for clustering heterogeneous data using rough set theory. *International Journal Rapid Manufacturing*, 1(2).
- Leibniz, G. W. (1989). *Discourse on Metaphysics*.
- Lenič, M., Povalej, P., & Kokol, P. (2005). Impact of Purity Measures on Knowledge Extraction in Decision Trees. In *Foundations and Novel Approaches in Data Mining*, May 2016, pp. 229–242. Berlin/Heidelberg: Springer-Verlag.
- Leung, Y., Fischer, M. M., Wu, W. Z., & Mi, J. S. (2008). A rough set approach for the discovery of classification rules in interval-valued information systems. *International Journal of Approximate Reasoning*, 47(2), 233–246.
- Li, J., Bioucas-Dias, J. M., Plaza, A., & Liu, L. (2016). Robust Collaborative Nonnegative Matrix Factorization for Hyperspectral Unmixing. *IEEE Transactions on Geoscience and Remote Sensing*, 9(9), 4267–4279.
- Li, L., Yang, J., Zhao, K., Xu, Y., Zhang, H., & Fan, Z. (2014). Graph Regularized Non-negative Matrix Factorization By Maximizing Correntropy. *JOURNAL OF COMPUTERS*, 9(11), 2570–2579.
- Li, T., & Ogihara, M. (2004). Entropy-Based Criterion in Categorical Clustering. In *Proceedings of the 21st International Conference on Machine Learning, Banff, Canada*.
- Liao, S. H., Chu, P. H., & Hsiao, P. Y. (2012). Data mining techniques and applications - A decade review from 2000 to 2011. *Expert Systems with Applications*, 39(12), 11303–11311.
- Lichman, M. (2013). UCI machine learning repository.
- Lingras, P. (2002). Rough set clustering for Web mining. *2002 IEEE World Congress on Computational Intelligence. 2002 IEEE International Conference on Fuzzy Systems. FUZZ-IEEE'02. Proceedings (Cat. No.02CH37291)*, 2, 1039–1044.

- MacQueen, J. B. (1967). K means some methods for classification and analysis of multivariate observations. *5th Berkeley Symposium on Mathematical Statistics and Probability 1967*, 1(233), 281–297.
- Maqbool, O., & Babri, H. A. (2007). Hierarchical clustering for software architecture recovery. *IEEE TRANSACTIONS ON SOFTWARE ENGINEERING*, 33(11), 759–780.
- Mathieu, R. G., & Gibson, J. E. (1993). A methodology for large-scale R&D planning based on cluster analysis. *IEEE Transactions on Engineering Management*, 40(3), 283–292.
- Mazlack, L. J., He, A., & Zhu, Y. (2000). A Rough Set Approach in Choosing Partitioning Attributes. In *Proceedings of the ISCA 13th, International Conference, CAINE*. pp. 1–6.
- Michael N. Tuma, Sören W. Scholz, R. D. (2009). The Application Of Cluster Analysis In Marketing Research. *Business Quest*.
- Miyamoto, S., & Takumi, S. (2012). Hierarchical clustering using transitive closure and semi-supervised classification based on fuzzy rough approximation. In *2012 IEEE International Conference on Granular Computing*. IEEE. pp. 359–364.
- Mohebi, E., & Sap, M. (2009). Rough Set Based Clustering of the Self Organizing Map. *2009 First Asian Conference on Intelligent Information and Database Systems*, (1), 82–85.
- Mudambi, S. (2002). Branding importance in business-to-business markets. Three buyer clusters. *Industrial Marketing Management*, 31(6), 525–533.
- Naresh Kumar Nagwani, S. V. (2012). CLUBAS: An Algorithm and Java Based Tool for Software Bug Classification Using Bug Attributes Similarities. *Journal of Software Engineering and Applications*, 05(06), 436–447.
- Naseem, R., Maqbool, O., & Muhammad, S. (2010). An Improved Similarity Measure for Binary Features in Software Clustering. In *Second International Conference on Computational Intelligence, Modelling and Simulation*.
- Naseem, R., Maqbool, O., & Muhammad, S. (2013). Cooperative clustering for



- software modularization. *The Journal of Systems & Software*, 86(8), 2045–2062.
- Ngo, C. L., & Nguyen, H. S. (2004). A Tolerance Rough Set Approach. *Knowledge Discovery in Databases*, pp. 515–517.
- Norušis, M. (2011). Cluster Analysis. In *Statistical Procedures Companion*, pp. 361–391.
- P. Danziger (2015). Big O Notation. Tech. rep.
- Park, I. K., & Choi, G. S. (2015a). A variable-precision information-entropy rough set approach for job searching. *Information Systems*, 48, 279–288.
- Park, I.-k., & Choi, G.-s. (2015b). Rough set approach for clustering categorical data using information-theoretic dependency measure. *Information Systems*, 48, 289–295.
- Parmar, D., Wu, T., & Blackhurst, J. (2007). MMR: An algorithm for clustering categorical data using Rough Set Theory. *Data & Knowledge Engineering*, 63(3), 879–893.
- Pawlak, Z. (1991). *Rough Sets Theoretical Aspects of Reasoning about Data*. Kluwer Academic Publishers.
- Pawlak, Z. (1995). Vagueness and uncertainty: A Rough Set Perspective. *Computational Intelligence*, 11(2), 227–232.
- Pawlak, Z. (1996). Rough sets and data analysis. In *Proceedings of Asian Fuzzy Systems Symposium on Soft Computing in Intelligent Systems and Information Processing*. IEEE. pp. 1–6.
- Pawlak, Z., & Skowron, A. (2007). Rudiments of rough sets. *Information Sciences*, 177(1), 3–27.
- Pawlak et al., Z. (1995). Rough sets. *Communications of the ACM*, 38(11), 88–95.
- Peters, G. (2006). Some refinements of rough k-means clustering. *Pattern Recognition*, 39(8), 1481–1491.
- Prabha, K., & Visalakshi, N. (2014). Improved Particle Swarm Optimization Based K-Means Clustering. *Intelligent Computing Applications*.
- Purwitasari, D., Fatichah, C., Ariesianti, I., & Hayatin, N. (2015). K-medoids

- algorithm on Indonesian Twitter feeds for clustering trending issue as important terms in news summarization. *Proceedings of 2015 International Conference on Information and Communication Technology and Systems, ICTS 2015*, pp. 95–98.
- Qamar, U. (2013). A Rough-Set Feature Selection Model for Classification and Knowledge Discovery. *Proceedings of IEEE International Conference on Systems, Man, and Cybernetics*, pp. 788–793.
- Rahman, M. N. a., Lazim, Y. M., & Mohamed, F. (2011). Applying Rough Set Theory in Multimedia Data Classification. *International Journal on New Computer Architectures and Their Applications (IJNCAA)*, 1(3), 683–693.
- Ramani, G. (2013). Rough set with Effective Clustering Method. *International Journal of Advanced Research in Computer and Communication Engineering*, 2(2), 1163–1167.
- Reddy, H. V., Viswanadha Raju, S., & Agrawal, P. (2013). Data labeling method based on cluster purity using relative rough entropy for categorical data clustering. In *Proceedings of International Conference on Advances in Computing, Communications and Informatics (ICACCI)*. IEEE. pp. 500–506.
- Rissino, S., & Lambert-torres, G. (2009). Rough Set Theory Fundamental Concepts , Principals , Data Extraction , and Applications. In Julio Ponce and Adem Karahoca (Ed.) *Data Mining and Knowledge Discovery in Real Life Applications*, pp. 35–58. I-Tech, Vienna, Austria.
- Rousseeuw, P. J. (1987). Silhouettes: A graphical aid to the interpretation and validation of cluster analysis. *Journal of Computational and Applied Mathematics*, 20(C), 53–65.
- Sachdeva, S., & Kastore, B. (2014). Document Clustering : Similarity Measures. Tech. Rep. 11693, Indian Institute of Technology Kanpur.
- Senan, N., Ibrahim, R., Nawi, N. M., Tri, I., Yanto, R., & Herawan, T. (2011). Rough Set Approach for Attributes Selection of Traditional Malay Musical Instruments Sounds Classification 1. *International Journal of database theory and Application*, 4(3), 59–76.

- Shelly, D. R., Hardebeck, J. L., Ellsworth, W. L., & Hill, D. P. (2016). A new strategy for earthquake focal mechanisms using waveform-correlation-derived relative polarities and cluster analysis: Application to the 2014 Long Valley Caldera earthquake swarm. *Journal of Geophysical Research: Solid Earth*, 121(12), 8622–8641.
- Shuanhu et al., W. (2004). Cluster Analysis of Gene Expression Data Based on Self-Splitting and Merging Competitive Learning. *IEEE Transactions on Information Technology in Biomedicine*, 8(1), 5–15.
- Singh, S., Mayfield, C., Prabhakar, S., Shah, R., & Hambruch, S. (2007). Indexing uncertain categorical data. In *International Conference on Data Engineering*. pp. 616–625.
- Sripada, S. C. (2011). Comparison of Purity and Entropy of K-Means Clustering and Fuzzy C Means Clustering. *Indian Journal of Computer Science and Engineering*, 2(3), 343–346.
- Sun, B., Yao, H., Ji, R., Xu, P., Sun, X., & Yuan, K. (2010). Individual Home-Video Collecting Using a Co-clustering Method. *First International Conference on Pervasive Computing, Signal Processing and Applications*, pp. 1132–1135.
- Suraj, Z. (2004). An Introduction to Rough Set Theory and Its Applications. In *ICENCO2004*.
- Tan, P. N., Steinbach, M., & Kumar, V. (2006). *Introduction to Data Mining*. Addison-Wesley.
- Tripathy, B., & Ghosh, A. (2011a). SDR: An algorithm for clustering categorical data using rough set theory. *IEEE Recent Advances in Intelligent Computational Systems*, pp. 867–872.
- Tripathy, B., Goyal, A., Chowdhury, R., & Sourav, P. A. (2017). MMeMeR: An algorithm for clustering heterogeneous data using rough set theory. *International Journal of Intelligent Systems and Applications*, 8, 25–33.
- Tripathy, B. K., & Ghosh, A. (2011b). SSDR : An Algorithm for Clustering Categorical Data Using Rough Set Theory. *Advances in Applied Science Research*, 2(3), 314–326.

- Tripathy, B. K., Goyal, A., & Sourav, P. A. (2016). A comparative analysis of rough intuitionistic fuzzy k-mode algorithm for clustering categorical data. *Research Journal of Pharmaceutical, Biological and Chemical Sciences*, 7(5), 2787–2802.
- Voges, K. E., & Pope, N. K. L. (2012). Rough Clustering Using an Evolutionary Algorithm. *Proceeding of 45th Hawaii International Conference on System Sciences*, pp. 1138–1145.
- Voges, K. E., Pope, N. K. L., & Brown, M. R. (2002). Cluster Analysis of Marketing Data: A Comparison of K-Means, Rough Set, and Rough Genetic Approaches. In *Heuristics and Optimization for Knowledge Discovery*, pp. 208–216.
- Wang, W., Gao, W., Wang, C., & Li, J. (2013). An Improved Algorithm for CART Based on the Rough Set Theory. In *Proceedings of Fourth Global Congress on Intelligent Systems*, January 2002. Ieee. pp. 11–15.
- Wang, Y., Liu, P., Guo, H., Li, H., & Chen, X. (2010). Improved Hierarchical Clustering Algorithm for Software Architecture Recovery. In *International Conference on Intelligent Computing and Cognitive Informatics*. pp. 1–4.
- Warren Liao, T. (2005). Clustering of time series data - A survey. *Pattern Recognition*, 38(11), 1857–1874.
- Wong, K.-P., Feng, D., Meikle, S. R., & Fulham, M. J. (2000). Segmentation of dynamic PET images using cluster analysis. *IEEE Symposium on Nuclear Science*, 3, 126–130.
- Wu, J., Hassan, A. E., & Holt, R. C. (2005). Comparison of clustering algorithms in the context of software evolution. In *IEEE International Conference on Software Maintenance, ICSM*, vol. 2005. pp. 525–535.
- Wu, J., Xiong, H., & Chen, J. (2009). Adapting the right measures for K-means clustering. *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining - KDD '09*, p. 877.
- Wu, J. S., Lai, J. H., & Wang, C. D. (2011). A novel co-clustering method with intra-similarities. *Proceedings of IEEE International Conference on Data Mining, ICDM*, pp. 300–306.

- Xu, H. Q., Besant, C. B., & Ristic, M. (2003). System for enhancing supply chain agility through exception handling. *International Journal of Production Research*, 41(6), 1099–1114.
- Xu, R., & Wunsch, D. (2005). Survey of clustering algorithms. *IEEE Transactions on Neural Networks*, 16(3), 645–678.
- Xue, G.-R., Lin, C., Yang, Q., Xi, W., Zeng, H.-J., Yu, Y., & Chen, Z. (2005). Scalable collaborative filtering using cluster-based smoothing. *Proceedings of the 28th annual international ACM SIGIR conference on Research and development in information retrieval - SIGIR '05*, p. 114.
- Yanto, I., Herawan, T., & Deris, M. (2011). Data clustering using variable precision rough set. *Intelligent Data Analysis*, 15, 465–482.
- Yanto, I. T. R., Ismail, M. A., & Herawan, T. (2016). A modified Fuzzy k-Partition based on indiscernibility relation for categorical data clustering. *Engineering Applications of Artificial Intelligence*, 53, 41–52.
- Zaïane, O. R. (1999). (Chapter 1) Introduction to Data Mining. *Principles of Knowledge Discovery in Databases*, pp. 1–15.
- Zhang, L., Li, Y., Sun, C., & Nadee, W. (2013). Rough Set Based Approach to Text Classification. In *Proceedings of International Joint Conferences on Web Intelligence (WI) and Intelligent Agent Technologies (IAT)*. Ieee. pp. 245–252.
- Zhao, Y. (2001). Criterion functions for document clustering: Experiments and analysis. Tech. rep., Department of Computer Science, University of Minnesota.
- Zhao, Y., & KARYPIS, G. (2004). Empirical and Theoretical Comparisons of Selected Criterion Functions for Document Clustering . *Machine Learning*, 55, 311–331.
- Zhong, S., & Ghosh, J. (2005). Generative model-based document clustering: a comparative study. *Knowledge and Information Systems*, 8(3), 374–384.
- Zhou, Z., & Mu, L. (2016). Representative Virtual Machine Templates: An optimized virtual machine templates management mechanism for an Cloud system based on K-medoids Clustering. In *Proceedings of 35th Chinese Control Conference*

(CCC). IEEE. pp. 5243–5248.

Zivkovic, Z. (2004). Improved adaptive Gaussian mixture model for background subtraction. In *Proceedings of the 17th International Conference on Pattern Recognition, 2004. ICPR 2004.*, vol. 2. IEEE. pp. 28–31.

Zsidsisin, G. a., & Ellram, L. M. (2003). An agency theory investigation of supply risk management. *The Journal of Supply Chain Management*, 39(August), 15–27.

